



Europäisches Patentamt
European Patent Office
Office européen des brevets



(11) **EP 0 817 170 A2**

(12) **EUROPEAN PATENT APPLICATION**

(43) Date of publication:
07.01.1998 Bulletin 1998/02

(51) Int Cl.⁶: **G10L 5/06, G10L 7/08,
G10L 9/06, G10L 9/18**

(21) Application number: **97850086.6**

(22) Date of filing: **05.06.1997**

(84) Designated Contracting States:
**AT BE CH DE DK ES FI FR GB GR IE IT LI LU MC
NL PT SE**

(72) Inventors:
• **Sundberg, Erik**
111 31 Stockholm (SE)
• **Melin, Hakan**
178 31 Ekerö (SE)

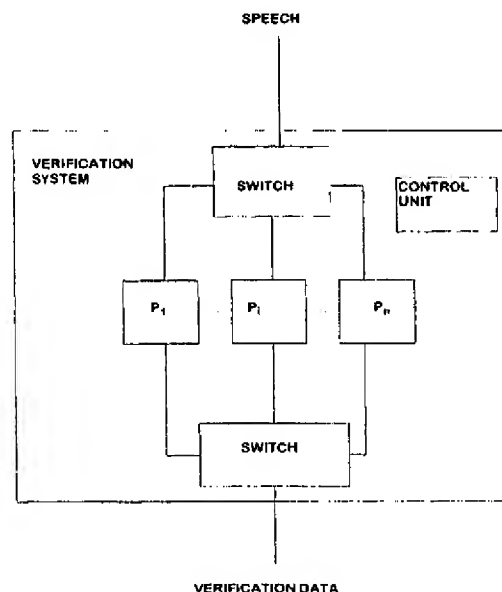
(30) Priority: **01.07.1996 SE 9602622**

(71) Applicant: **TELIA AB**
123 86 Farsta (SE)

(74) Representative: **Karlsson, Berne**
Telia Research AB,
Rudsjösterrassen 2
136 80 Haninge (SE)

(54) **Method and apparatus for adaption of models in e.g. speaker verification systems**

(57) The invention relates to a method and an arrangement for adapting models in speaker verification systems or similar systems using models based on data collected from a person during a certain time period. If a simple model is utilized a less reliable verification is obtained, but if, on the other hand, a more complex model is utilized, the problem is a long training time period. The present invention solves this problem by using a plurality of different models in the same speaker verification system. The verification is put into operation using the model requiring a smaller amount of speech data. During the use, more speech data is collected continuously. This material is used to further train either only the more complex model or both the simpler model already put into operation and the more complex model. At suitable intervals, a comparison is made of the performance capacities of the models. Once the more complex model yields a more reliable verification result it will take over the verification in the operating situation. The subsequent model unit may be put into operation either instantaneously or gradually, e.g. by using a weight function.



FIG

EP 0 817 170 A2

Description

FIELD OF THE INVENTION

The present invention relates to a method and an arrangement to adapt models in speaker verification systems or similar systems using models based on data collected from a person during a certain time period. The collected data may be related to the physiology, behavior, aging of the person etc. A related area is e.g. speaker adaptive speech recognition. In systems of this type, collected data is compared to a model for the verification of the identity of the speaker or recognition of the speech in order to control a course of events in a process or of a device. For the model to be able to perform its task it has to be trained with speech data. Simpler models require less training but provide a less reliable result, while more complex models require longer training and provide a more reliable result of the verification.

The invention may be applied in all speaker verification systems that are to be used at a plurality of occasions, that is speech of the same person is to be verified at repeated occasions. As is known, speaker verification systems are utilized in order to protect information or economic values. The invention is an alternative to the approach of using PIN codes in order to identify a user. The voice recording as such may be effected either directly at the equipment where the verification is performed or is transmitted by various media. The medium may be telephony or other telecommunication media.

STATE OF THE ART

In the prior art speaker verification systems use has been made of only one model with the special problems associated with the model. Thus, if a simple model has been used, a less reliable verification has been obtained. On the other hand, if a more complex model is used, the problem is a long term training period.

The present invention solves this problem by utilizing a plurality of different models in the same speaker verification system. The verification is put into operation with the model requiring the less amount of speech data. During the use, more speech data is continuously collected. This material is used to further train either only the more complex model or both the simpler model already in operation and the more complex model. At suitable points of time, comparisons are made of the performances of the models. When the more complex model provides a more reliable verification result, it will take over the verification in the operation situation.

It is recognized that because of the invention a speaker verification system is obtained that is readily put into operation but eventually providing increasingly reliable verification results. The invention enables a use of the advantages of different models at the same time as the effect of their respective disadvantages is mini-

mized. Without using this technology, one has to choose a model with its associated advantages and disadvantages at the start of a speaker verification system. By shifting models it is achieved that the system dynamically adapts to the available amount of speech data. This means a great advantage over the prior art.

SUMMARY OF THE INVENTION

Thus, the present invention provides a method for adapting a model in e.g. speaker verification, comprising model units for receiving and evaluating speech. According to the invention, speech data is collected and a first model unit is put into operation while a subsequent model unit is trained with speech data being collected during the operation of the first model unit. The performances of the model units are tested and evaluated and a subsequent model unit is put into operation when the performance thereof has reached a suitable level.

The subsequent model unit may be put into operation either instantaneously or gradually, e.g. by using a weight function.

The invention also relates to an arrangement for performing the method.

The invention is defined in detail in the accompanying claims.

BRIEF DESCRIPTION OF THE DRAWING

The invention will be described in detail below with reference to the attached drawing, wherein the only Figure is a schematic illustration of an embodiment of the invention.

DETAILED DESCRIPTION OF A PREFERRED EMBODIMENT

In speaker verification systems, systems for the automatic verification of the identity of a speaker, the amount of speech data that has to be collected from the user is a determining limitation of the use. Complex speaker models requiring a large amount of collected speech data may be expected to provide a better result than models requiring a small amount of training material. However, for a small amount of training material the more complex model may yield an inferior result than the simpler model.

Complex models having many parameters have better performance than simpler models once the parameters of the model has been estimated correctly. However, for a correct estimation of the parameters a large amount of training data is required. In the case where the training data of a model is provided by a customer, the amount of training data is a factor of inconvenience for the customer. Poor performance within a model will also lead to system errors, being another factor of inconvenience for the customer. A problem that is solved by the present invention is to find model topolo-

gies having good performance with a minimum of training data.

The solution of the problem proposed herewith of both maximizing the performance of the model and minimizing the requirement of training data is to use a model system having a dynamic topology. The model has a combination of model units or parts having varying degree of complexity. The effective topology of the model is changed, such that for a given amount of training data the optimum topology is used, based on the given model unit. By using this technique, the effective complexity of the model will grow with the available amount of training data.

In the beginning of the service life of the model, the simplest model units are used, requiring only a small amount of data for a reliable estimation of its parameters. As the amount of available data grows, the more complex parts can be trained successively.

Once the parameters of the more complex unit have been estimated in a reliable way, the performance thereof is probably better than that of the simpler unit and the topology of the model may be changed in favour of the complex unit.

In the single Figure a speaker verification system in accordance with the present invention is illustrated schematically. The system comprises a control unit controlling two switches and a number of model units P_1 - P_n . On the one hand, the system receives speech or speech data and supplies verification data as the output signal.

The various model units P_1 - P_n of the speaker model have different requirements of training data. A model unit P_i should only be used for verification when it has received sufficient training data. The units requiring a smaller amount of data will be put into operation earlier, while the more demanding units will not be used until a longer training period has elapsed. In this way, the performance of the speaker model may be enhanced towards its full capacity. During the growth period the model may still be used for verification by using the simpler model units of the speaker model.

The simpler parts may be taken out of service as the more complex units achieve better performance.

The shift to newer models may be effected over several generations, so that more and more advanced models requiring more speech data continuously is put into operation. In this way, the speaker verification system may be upgraded without being put out of operation. In addition, it is contemplated that each model consists of several submodels weighted together in various ways to define a model.

When the speaker verification system is put into operation the very first time it requires a short training period to train the simplest model unit. The simplest model unit can be trained from a speaker independent template. Thereafter, the system is put into operation with increasing performance in accordance with what is stated above.

Each unit of the speaker model hierarchy will need to store information relating to how well trained it is. This information may be provided either by the model unit itself or by some performance testing method. In the former case, the information is called training level while, in the latter case, the information is called performance level. The training level is based on an assumed a priori knowledge about how much training data is needed by the unit. The difference between the two kinds of information is that the performance level is based on some evaluation of test data (a data base run), while the training level is based on stored information about used training data. The performance level may be based on comparisons with other units of the speaker model and even other speaker models.

Thresholds for the training level and the performance level must be provided and stored in the control unit. In the former case, the threshold is based on previously made assumptions. For the latter, it should be possible to base the value of the threshold on a criterion of the performance demand.

In order to enable use of a performance level based on data base simulation, it is necessary to include management of such a data base. The speaker model should also be able to state a value of its total training level or performance level. This value may be used by an application to estimate the significance level of a decision taken by the verification system.

The performance of the model units is tested at suitable intervals in order to check if they should be operative or not. This may be effected cyclically or on a special command.

The invention has been described with reference to a speaker verification system but, as is mentioned above, the invention may equally be applied in other systems using models based on data collected from a person under a certain time period, e.g. speaker adaptive speech recognition systems. The invention is only limited by the claims below.

Claims

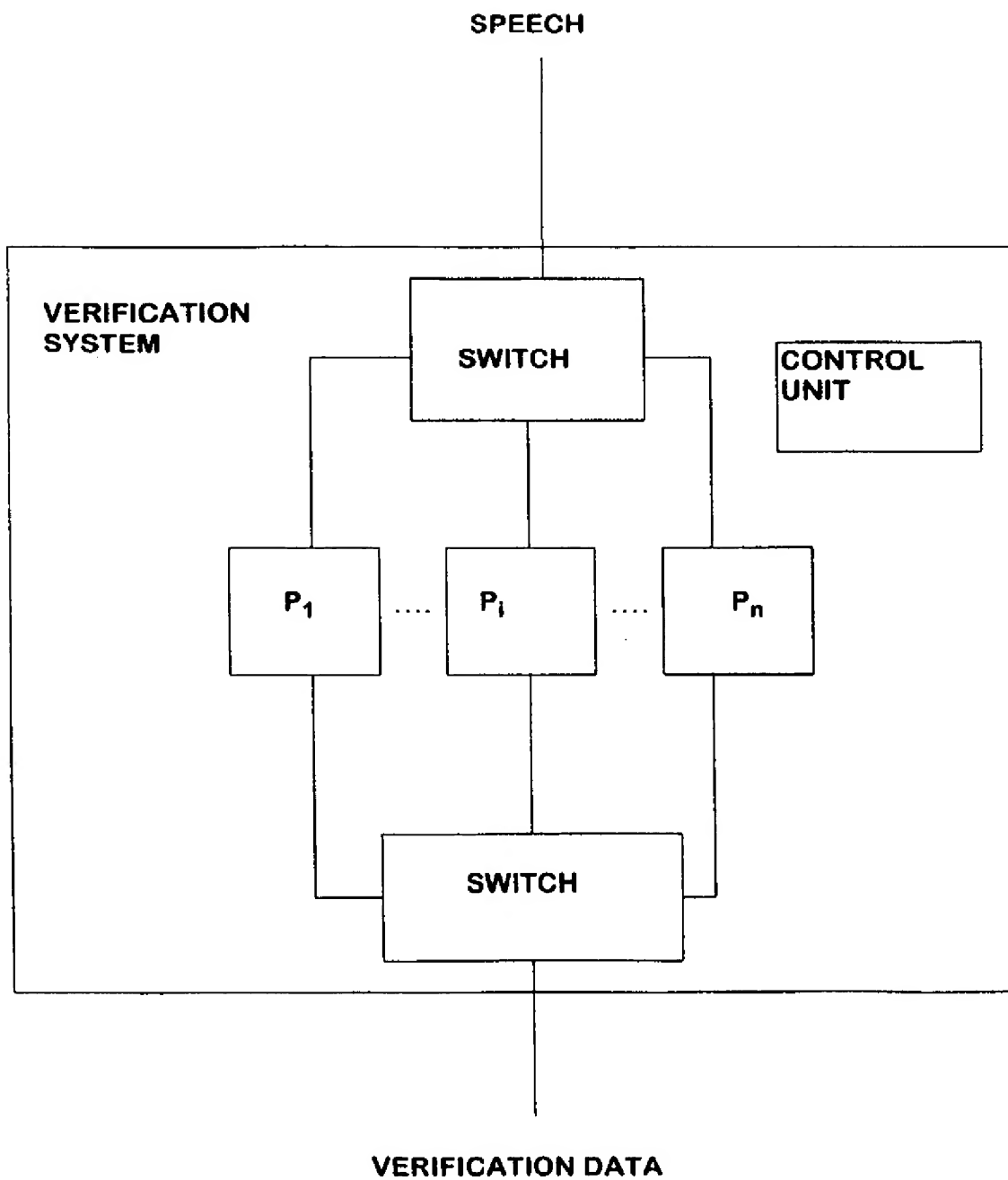
1. Method for adapting a model in e.g. speaker verification systems, comprising model units for receiving and evaluating speech, characterized by collecting speech data and putting a first model unit (P_1) into operation, training a subsequent model unit (P_n) with speech data collected during the operation of the first model unit, testing and evaluating the performance capacities of the model units and putting the subsequent model unit into operation once the performance capacity thereof has reached a suitable level.
2. Method in accordance with claim 1, characterized by putting the subsequent model unit into operation instantaneously once the performance capacity

thereof exceeds a predetermined threshold.

3. Method in accordance with claim 1, characterized by putting the subsequent model unit into operation gradually, once the performance capacity thereof exceeds a respective threshold of a plurality of predetermined thresholds. 5
4. Method in accordance with claim 3, characterized by putting the subsequent model unit into operation gradually by weighting the various model units with a variable weight function. 10
5. Method in accordance with anyone of the previous claims, characterized by connecting a new model unit as a subsequent model. 15
6. Method in accordance with anyone of the previous claims, characterized by training all model units with collected speech data. 20
7. Method in accordance with anyone of claims 1 to 6, characterized by training all model units except the respective operative model units with collected speech data. 25
8. Arrangement for adapting a model in e.g. speaker verification systems, comprising model units for receiving and evaluating speech and a control unit, characterized by a first switch for directing speech data to the various model units (P_1 - P_n), a second switch for directing verification data from the various model units (P_1 - P_n), said switches being controlled by the control unit such that the model units collect speech data and that a first model unit (P_1) is put into operation, a subsequent model unit (P_n) is trained with speech data collected during the operation of the first model unit, that the performance capacity of the model units are tested and evaluated and that the subsequent model unit is put into operation once the performance capacity thereof has reached a suitable level. 30 35 40
9. Arrangement in accordance with claim 8, characterized in that a predetermined threshold is stored in the control unit in order to put the subsequent model unit into operation instantaneously, once the performance capacity thereof exceeds the predetermined threshold. 45 50
10. Arrangement in accordance with claim 8, characterized in that a plurality of predetermined thresholds are stored in the control unit in order to put the subsequent model unit into operation gradually, once the performance capacity thereof exceeds a respective threshold of the predetermined plurality of thresholds. 55

11. Arrangement in accordance with claim 10, characterized in that the control unit comprises a variable weight function to put the subsequent model unit into operation gradually by weighting the different models with the weight function.

12. Arrangement in accordance with anyone of claims 8 to 11, characterized in that a model unit consists of submodels.



FIG